

The Drosophila Gene Disruption Project: Progress Using Transposons With Distinctive Site-Specificities

Hugo J. Bellen*, Robert W. Levis†, Yuchun He*, Joseph W. Carlson§, Martha Evans-Holm§, Eunkyung Bae‡1, Jaeseob Kim‡2, Athanasios Metaxakis**3, Charalambos Savakis**4, Karen L. Schulze*, Roger A. Hoskins§ and Allan C. Spradling†

*Howard Hughes Medical Institute
Department of Molecular and Human Genetics
Program in Developmental Biology
Baylor College of Medicine
Houston, TX 77030

†Howard Hughes Medical Institute Research Laboratories
Department of Embryology
Carnegie Institution for Science
Baltimore, MD 21218

§Lawrence Berkeley National Laboratory
Life Sciences Division
Berkeley, CA 94720

‡Aprogen Inc.
Seoul, KOREA

****Institute of Molecular Biology and Biotechnology**
Foundation for Research and Technology
Heraklion 71110
Crete, GREECE

¹Present address: Schnell Biopharmaceuticals, Inc. Sincheon-Dong 11-10, Songpa-Gu, 138-240, Seoul, KOREA email: eunkyung80@schnell.kr

²Full address: Aprogen Inc., Sangdaewon-Dong 442-2, Joongwon-Gu, Sungnam-Shi, Kyunggi-Do, 462-807, Seoul, KOREA email: jaeseob@aprogen.com

³Present address: Max Planck Institute for Biology of Ageing, D-50931 Cologne, GERMANY

⁴Present address: BSRC Alexander Fleming, P.O. Box 74145, 16602 Varkiza, Greece

Running Head: Drosophila gene disruption project

Keywords: *P*-element, *Minos*, insertion, mutation

Corresponding author:

Allan Spradling
Department of Embryology
Carnegie Institution
Baltimore, MD 21218
Email: spradling@ciwemb.edu

ABSTRACT

The *Drosophila* Gene Disruption Project has created a public collection of mutant strains containing single transposon insertions associated with different genes. These strains often disrupt gene function directly, allow production of new alleles, and have many other applications for analyzing gene function. Here we describe the addition of about 7,600 new strains, which were selected from more than 140,000 additional *P* or *piggyBac* element integrations and 12,500 newly generated insertions of the *Minos* transposon. These additions nearly double the size of the collection and increase the number of tagged genes to at least 9,440, approximately two-thirds of all annotated protein-coding genes. We also compare the site-specificity of the three major transposons used in the project. All three elements insert only rarely within many Polycomb-regulated regions, a property that may contribute to the origin of “transposon-free regions” in metazoan genomes. Within other genomic regions, *Minos* transposes essentially at random whereas *P* or *piggyBac* elements display distinctive hotspots and coldspots. *P* elements, as previously shown, have a strong preference for promoters. In contrast, *piggyBac* site selectivity suggests that it has evolved to reduce deleterious and increase adaptive changes in host gene expression. The propensity of *Minos* to integrate broadly makes possible a hybrid finishing strategy for the project that will bring >95% of *Drosophila* genes under experimental control within their native genomic contexts.

FUNDING: Department of Energy under Contract No. DE-AC02-05CH11231.

DISCLAIMER

This document was prepared as an account of work sponsored by the United States Government. While this document is believed to contain correct information, neither the United States Government nor any agency thereof, nor the Regents of the University of California, nor any of their employees, makes any warranty, express or implied, or assumes any legal responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by its trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof, or the Regents of the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof or the Regents of the University of California.

INTRODUCTION

Drosophila has served as an important model organism for over 100 years, in large part, because of the wealth of mutants available and the ease with which they can be manipulated experimentally. Mutagenesis using single insertions of an engineered transposon offers many advantages for analyzing gene regulation and function (COOLEY et al. 1988; BELLEN et al. 1989; BIER et al. 1989). The insertions frequently interfere directly with gene function, and can also be re-mobilized to generate additional useful mutations in the genomic region where they reside through the processes of local jumping or imprecise excision. By incorporating useful internal sequences, transposons can be used to report or manipulate gene expression, sense chromatin structure, or function as sites for site-specific recombination.

The *Drosophila* Gene Disruption Project (GDP) was established in 1991 to bring the advantages of this method to the research community by generating transposon mutations in most *Drosophila* genes. During phase 1 of the project, we characterized insertions causing recessive phenotypes (SPRADLING et al. 1999). The availability of an annotated genome sequence (ADAMS et al. 2000; MISRA et al. 2002) enabled phase 2, where insertions were associated with predicted genes based solely on their genomic location. By 2004 about 40% of known *Drosophila* genes had one or more associated GDP insertion alleles (BELLEN et al. 2004). Several large collections of insertion lines were independently generated as well, further increasing the potential gene coverage (THIBAUT et al. 2004; KIM et al. 2010).

Several different approaches may help further increase the number of disrupted genes. Transposable elements differ in their target site-specificity (THIBAUT et al. 2004; BELLEN et al. 2004); hence generating insertions using new transposons might provide greater efficiency than continued mutagenesis using *P* and *piggyBac* elements. The *Minos* transposon, a mariner family

member, is a particularly attractive candidate (METAXAKIS et al. 2005). The site-specific recombinase from phage ϕ C31 provides the ability to efficiently integrate even large DNAs into genomic *attP* target sites (GROTH et al. 2004; BATEMAN et al. 2006; VENKEN et al 2006; 2009). Including a ϕ C31 *attP* site in the elements used for mutagenesis would offer many advantages for genomic manipulation, including increased mutagenicity (GROTH et al. 2004; BATEMAN et al. 2006; VENKEN et al. 2006). Previous studies of the capabilities of integrated *attP*-containing transposons illustrate their exceptional utility (VENKEN AND BELLEN, 2007; VENKEN et al. 2009).

Transposon site-specificity represents a critically important factor in determining the optimum strategy for completing the GDP project. The size and quality of the data collected by the GDP provide a special opportunity to characterize the insertional preference of specific transposons in detail. It is well established that some transposons hit certain sites, “hotspots,” much more frequently than expected by chance, while other regions, “coldspots,” are avoided. *P* elements frequently insert near promoters, an advantage for mutagenesis and mis-expression screening, but also preferentially target hotspots (SPRADLING et al. 1995; LIAO et al. 2000; BELLEN et al. 2004). In addition, a significant fraction of *Drosophila* genes, including many clustered and tissue-specific genes, appear almost refractory to disruption by *P* insertion (BELLEN et al. 2004). *piggyBac* elements also target hotspots, but show less regional and promoter bias (THIBAUT et al. 2004; BELLEN et al. 2004). However, *piggyBac* elements do not excise imprecisely to create local deletions, a significant disadvantage compared to *P* elements.

Here we summarize the status of the GDP collection at the completion of phase 2. We have added *P*-element, *piggyBac* and *Minos* insertions to the publicly available GDP collection to provide genetic access to at least 9,440 genes. In addition to expanding this resource of

mutants for researchers, our studies also provide new insights into transposon site-selectivity and document an influence of chromatin structure. We show that, because of very low site-specificity, it should be feasible to tile the *Drosophila* genome with *Minos* insertions that would facilitate the site-directed mutagenesis of almost all *Drosophila* genes and functional elements by homologous recombination.

MATERIALS AND METHODS

The EY collection: The construction of the *P*-element based EY transposon (*P{EPgy2}*, Table 1), the generation of 10,310 insertion lines, and the mapping of their insertion sites have been described (BELLEN et al. 2004). Using the same methods we generated an additional 11,830 EY lines (strain names EY10505-EY16964 and EY18301-EY23670) and mapped 9,585 insertions to unique sites on the reference *Drosophila* genomic sequence (Release 5, <http://www.fruitfly.org>). This brought the total number of unselected EY transpositions generated and uniquely localized in the genome to 18,214. The new EY insertions that we selected for the GDP collection were balanced and their insertion sites were verified by resequencing of flanking DNA as described (BELLEN et al. 2004).

The Exelixis collection: The generation and properties of 26,540 *P* and *piggyBac* insertion lines that were mapped by Exelixis, Inc. to unique sites in the *Drosophila* reference genomic sequence (Release 2) have been described (THIBAUT et al. 2004). These lines probably do not represent a completely random collection of insertions, because some lines disrupting major hotspots appear to have been culled by Exelixis. However, we found many cases where at least two lines bearing identical *piggyBac* insertion sites had been retained, suggesting that such culling was limited or incomplete. Most of these stocks, as well as insertion site data, were generously made available to the GDP so that the most useful lines could be

distributed publicly. Approximately 400 base pairs (bp) of the genomic reference sequence surrounding the insertion site(s) in these lines along with a coordinate or range of coordinates denoting the insertion site were reported (THIBAUT et al. 2004). We selected approximately 2,100 lines from the Exelixis collection for distribution by the Bloomington Drosophila Stock Center (BDSC), based on the insertion site coordinates reported by Exelixis. Exelixis subsequently provided us with 52,183 flanking sequence reads derived from 22,144 strains with associated phred quality scores (EWING AND GREEN, 1998). In April 2005, 24,678 of these flanking sequence reads were submitted to GenBank by the GDP. We subsequently realigned the flanking sequences to the Drosophila reference genomic sequence (Release 5) based on more stringent criteria using our standard pipeline and mapped 16,073 insertions to unique sites. While there was usually close agreement, insertion site coordinates deduced by the GDP and Exelixis sometimes varied by several hundred base pairs, and 535 strains lacked any sequence reads. Some strains had multiple sequence reads from one or both flanks and these sometimes mapped to different sites. After changes due to the reanalyzed sequence flanks, updated annotation, strain losses, and line substitutions, 1,859 Exelixis lines are currently part of the GDP collection at the BDSC, while 357 Exelixis GDP lines are maintained at Harvard Medical School (<https://drosophila.med.harvard.edu/>) (Table 2).

The MB collection: To generate new insertions of a *Minos* element, we used the *Mi{ET1}* element described in (METAXAKIS *et al.* 2005) (Table 1). It contains the *Minos* 255 bp inverted repeats and a minimal *hsp70* promoter upstream of the *GAL4* gene and may function as an enhancer detector/trap (hence “ET”) if inserted in the appropriate location. The GFP gene, driven in the eye and brain of adults and larvae by the 3xP3 promoter (HORN et al., 2000), is the marker used for selection. The stocks were generated and balanced in the *w¹¹¹⁸* isogenic

background described in (RYDER *et al.* 2004). The *Minos Mi{ET1}* mutator (FlyBase ID FBtp0021506; referred to as MiET1 by Metaxakis *et al.* 2005), which we refer to as the MB element, was inserted on a *TM3, Sb Ser* balancer chromosome. The starting site of the mutator was mapped by flanking sequence (GenBank accession ET202027) to a site corresponding to coordinate 3L:12580323 of the *D. melanogaster* reference genomic sequence. The MB mutator was mobilized using a transgenic source of transposase under the control of a heat-shock promoter (*P{hsILMiT}*, FlyBase ID FBtp0021508, referred to as PhsILMiT by METAXAKIS *et al.* (2005) inserted on a second chromosome balancer (*P{hsILMiT}2.4*; FlyBase ID FBti0073645).

We generated 12,426 strains containing new insertions of the MB transposon (nearly always single insertions) and mapped 10,781 insertions from 10,630 strains to a unique site in the genome. Lines that were selected for the GDP collection were balanced and their insertion sites verified by resequencing before delivery to the BDSC. Sequences flanking MB insertions were determined by inverse PCR and DNA sequencing, as described in (BELLEN *et al.* 2004) with the following modifications. Genomic DNA was digested with Hpa II; 5' flanks were amplified with the primers MI.5.F (CAAAAGCAACTAATGTAACGG) and MI.5.R (TTGCTCTTCTTGAGATTAAGGTA) at an annealing temperature of 50°; 3' flanks were amplified with MI.3.F (ATGATAGTAAATCACATTACG) and MI.3.R (CAATAATTTAATTAATTTCCC) at an annealing temperature of 50°; 5' and 3' flanks were sequenced with MI.seq (TTTCGTCGTGAAGAGAAT). A detailed protocol is available on the GDP website (<http://flypush.imgen.bcm.tmc.edu/pscreen/>). Insertion-bearing chromosomes were balanced using *P{RS3}l(1)CB-6411-3'*, *w¹¹¹⁸/FM7h* (X chromosome), *w¹¹¹⁸/Dp(1;Y)y⁺*; *noc^{Scd}/SM6a* (2nd chromosome), and *w¹¹¹⁸/Dp(1;Y)y⁺*; *TM2/TM6C, Sb¹* (3rd chromosome), which are all in the “iso31” isogenic background (RYDER *et al.* 2004) and were obtained from the

BDSC. A Meme analysis failed to uncover any significant target sequence preference beyond the requirement for “TA.”

The GenExel (Aprogen) collection: GenExel, Inc., now Aprogen, Inc., generated a very large collection of lines bearing insertions of the *P*-element construct *P{EP}* (Rorth 1996; FlyBase ID FBtp0001317) at the Korean Advanced Institute of Science and Technology (KAIST); (Table 1, “G”; see <http://www.oxfordjournals.org/nar/database/summary/677>; KIM et al. 2010). Initially, approximately 27,000 lines were selected from a starting set of about 100,000 transpositions by requiring a minimum spacing of 200 bp between insertions in order to prune out lines with insertions in transposon hotspots. Most insertions were not balanced. Sequence coordinates for 24,789 insertions were provided to the GDP. GenExel subsequently sent us 1,685 strains that we had identified as candidates. After balancing the insertions and sequencing their flanks, 1,136 lines were added to the GDP collection at the BDSC.

The Max Planck/EMBL/DeveloGen collection: We received lines from a collection of *P* element insertions generated by researchers at Max Planck Göttingen, the EMBL labs at Heidelberg and DeveloGen AG (STAUDT et al. 2005). The lines comprising this collection are indicated by the prefixes HP or DP (Table 1). Insertion site information was provided, and lines hitting novel genes were identified for transfer directly from Max Plank to the BDSC.

Other collections: The Göttingen collections of insertions on the X chromosome (PETER et al. 2002; BEINERT et al. 2004) were screened. The elements comprising this collection are designated by the prefix G0 or GG (Table 1). Candidates from the *P{lacW}* insertion collection on FRT-bearing chromosomes described by OH et al. (2005) were resequenced and screened. The elements comprising this collection are designated by the prefix SH (Table 1).

Strain selection: Strains were selected for inclusion essentially as described previously

(BELLEN et al. 2004). The GDP employs a strategy of continuous library improvement, both by adding new lines, and by replacing/upgrading existing lines with better ones. Briefly, each new candidate insertion from the screens described is compared with the *Drosophila* genome annotation, as well as with the insertion sites of all existing GDP collection strains within the gene region in question. Based on the best judgment of an expert annotator, lines can be retained for several reasons. Of highest priority are lines likely to disrupt any gene lacking a current GDP insertion. Because the annotated 5' end of many gene models may be truncated relative to the true 5' end, insertions located within 500 bp of the annotated 5' end or anywhere within the transcribed region are selected. In addition, a second insertion in a gene is saved if it is located in a distinct promoter, disrupts another transcript isoform or provides another unique genetic property. The continued presence of unannotated protein-coding and RNA genes, and genetic regulatory elements, especially in annotations prior to modENCODE (ROY et al. 2010), provides the final reason for selecting lines. Since *P* elements show a strong preference for promoters, *P*-element insertions located 2 kb or more from the nearest annotated promoter or existing insertion are also retained. Similarly, a small number of *piggyBac* or MB lines that had insertions within regions more than 10 kb distant from any existing insertion have also been kept for use in genetically manipulating the surrounding genome. Many insertions thought initially to be within intergenic regions have subsequently been mapped to genes as the annotation improved. Many such lines have been used to functionally characterize novel genes, promoters, piRNA clusters and small RNA genes (for example BRENNECKE et al. 2003; GODFREY et al. 2006; BRENNECKE et al. 2007).

The GDP recognizes that lines added to the collection based on the above criteria are not equally valuable. Hence, lines whose value is less certain are subject to replacement. For

example, lines mapping upstream from annotated transcription units are replaced when lines became available whose insertions are located within the unit. Strains containing two insertions on the same chromosome are retained if one is located within a novel gene. However, such lines are also replaced as soon as a single copy insertion in the gene becomes available. Other reasons for line replacement are restraints on distribution. Some donated collections cannot be distributed to for-profit corporations. These lines are subject to replacement whenever an equivalent line without such conditions becomes available.

Data handing and access: Genomic sequences flanking the *P-element* and *piggyBac* transposon insertions were determined as described in BELLEN et al. (2004); sequences flanking MB insertions were determined as described above. The analysis and alignment of all flanking sequences were as described in BELLEN et al. (2004). The genome sequence coordinates given here are based on the Release 5 reference genome sequence. We consider an insertion to hit a gene if the insertion site is within the annotated transcription unit of the gene or within 500 bp upstream of the 5' end, based on the FlyBase gene annotation release FB2009_10.

The GDP website (<http://flypush.imgen.bcm.tmc.edu/pscreen/>) has a searchable database of strains that are part of the GDP collection at the BDSC, as well as those that have been selected to be added to the collection and are in the process of being balanced and rechecked. Data presented are the transposon construct, line name, genomic insertion site, inferred cytogenetic map location, associated gene, FlyBase annotation reference, and BDSC stock number.

Project data are sent to FlyBase (<http://flybase.bio.indiana.edu/>) and GenBank (<http://www.ncbi.nlm.nih.gov/>) before the lines are transferred to the BDSC for public distribution (<http://flystocks.bio.indiana.edu/>). Insertion data are displayed using the UCSC

genome browser (FUJITA *et al.* 2010). Custom tracks for this display are available from the GDP website. Complete insertion information on EY, MB and Exelixis *piggyBac* insertions that were analyzed in this study for site-specificity is available on the GDP website.

RESULTS

New *P* element and *piggyBac* insertion lines: Previous efforts generated a GDP collection consisting of 7,140 lines bearing *P* element or *piggyBac* insertions that provided access to 5,362 genes (BELLEN *et al.* 2004). One approach to further expanding the collection is simply to screen more lines containing unselected insertions of these elements. To this end, 11,830 new insertions of the EY element, a modified *P* transposon that can be used to misexpress endogenous genes adjacent to its insertion site (Table 1), were generated. In addition, two large collections of insertion strains were donated to the project. Exelixis, Inc. provided site coordinates for 6,194 *P* element and 18,668 *piggyBac* insertion lines. The structure of the *P{XP}*, *PBac{PB}*, *PBac{RB}* and *PBac{WH}* transposons used to construct these lines (THIBAUT *et al.* 2004) is shown in Table 1. GenExel, Inc. (currently, Aprogen Inc.) generously made available sequence coordinates from about 24,789 *P{EP}* element insertions (Table 1) that they selected from a starting collection of approximately 100,000 lines. Several other groups of investigators provided coordinates for smaller but significant collections (see METHODS).

The insertion sites in all the new lines, which include the full genetic diversity generated by more than 140,000 *P* and *piggyBac* element transpositions, were screened against the *Drosophila* genome annotation to identify lines that would expand the genetic diversity of the GDP collection. Overall, 5,002 *P* or *piggyBac* lines were added to the collection because their insertions were located in novel genes (3,439), in putative regulatory regions, or because they were more likely than a currently existing allele to strongly disrupt gene function (Table 2, see

METHODS for further details).

Generation of *Minos* insertion lines: Our results illustrate how random forward mutagenesis becomes increasingly inefficient as saturation is approached. About 50,000 *P* and *piggyBac* lines were required to identify insertions associated with the first 5,362 genes (BELLEN *et al.* 2004). Subsequently, our screening of nearly three times as many insertions yielded only 0.66 times as many new genes, highlighting the fact that *P* and *piggyBac* insertion sites were becoming saturated. Indeed, less than 2% of newly generated EY insertions near the end of the screen disrupted genes not previously represented in the collection.

To continue improving the GDP collection and to further investigate the options for finishing the project, a screen was carried out using *Minos*, a mariner family transposon unrelated to either *P* or *piggyBac*. A previous study (METAXAKIS *et al.* 2005) suggested that *Minos* integrates into the *Drosophila melanogaster* genome with little site-specificity. However, this conclusion was based on a small sample of about 100 insertions. To exploit the properties of this element and to measure its behavior more accurately, we carried out a large screen to generate new insertions using the *Minos*-based *Mi{ET1}* element (METAXAKIS *et al.* 2005; see Table 1). We refer to these as MB lines. Of the 12,426 MB lines with independent transpositions that were generated and sequenced, we recovered flanking sequence that could be unambiguously localized to a unique site in the genome from 10,630 lines (86%).

We added 2,658 of the MB lines to the GDP collection (Table 2). Although lines were saved for a variety of reasons, 1,155 of the MB lines hit genes new to the GDP collection, bringing the total number of disrupted genes to 9,440, which is about two-thirds of currently annotated *Drosophila* protein-coding genes (TWEEDIE *et al.* 2009). Thus, since the last report (BELLEN *et al.* 2004), the number of lines in the GDP collection has approximately doubled, and

the number of disrupted genes has increased by 77% (Figure 1).

Comparing the insertional specificities of *P*, *piggyBac* and *Minos* elements: The high efficiency of the MB screen in generating useful new insertions provided further evidence that significant differences exist in the insertional specificities of *P*, *piggyBac* and *Minos* elements. To further investigate whether to continue with forward *Minos*, we analyzed the site-specificities of MB, EY and *piggyBac* elements in detail. We used information from 18,214 EY insertions, 12,244 Exelixis *piggyBac* insertions that upon reanalysis by GDP were unambiguously mapped to unique sites, and 10,458 MB insertions. Both the EY and the MB screens incorporate data on all transpositions outside the chromosome bearing the starting insertion. In contrast, some redundant or nearly redundant *piggyBac* insertions may have been culled from the data sent by Exelixis (see METHODS). However, removal of lines with similar insertion sites would only serve to increase the apparent randomness of *piggyBac* insertion. In addition, the methods we used in generating and analyzing these data minimize problems caused by insertions within repetitive sequences or within heterochromatic regions that suppress marker gene expression (see DISCUSSION).

The MB screen showed one anomaly with the potential to skew our analysis. A general scan of the insertion distribution revealed the presence of a single large MB hotspot in chromosome 3L at 12.583 Mb, which corresponds to the site of the starting element located on a balancer chromosome homolog (Figure 2A). Such “homolog hotspots” have been observed previously in some, but not most *P*-element screens (TOWER and KURAPATI 1994; BELLEN et al. 2004). Approximately 310 of the 10,458 insertions were located within 300 kb of the starting site in a peaked distribution (Figure 2B). A similar distribution of new insertions arising near the original insertion on the starting chromosome has previously been observed when transposons

were experimentally re-mobilized, a phenomenon known as “local transposition.” However, homolog hotspots differ in that they result from hopping to nearby sites on the homolog, rather than the starting chromosome itself. No homolog hotspot was observed in the EY screen. Since this hotspot does not reflect the intrinsic site-specificity of *Minos* elements, these 310 lines were not used in analyzing site-specificity. However, these observations do provide evidence that *Minos* elements can undergo high frequency local transposition.

The insertional specificities of *P*, *piggyBac* and *Minos* elements differ: To visualize differences in transposition specificity we divided the 117 Mb “core” genome (including all euchromatin and some telomeric and pericentric heterochromatin) into regular 10 kb intervals and determined how many times each interval was hit by MB, *piggyBac* or EY insertions. To facilitate comparison, the same number of insertions was scored in each case (selected in numerical order by strain name), and this was set equal to the number of intervals (defined as $\lambda = 1$). Because the number of insertions on each arm varied, each arm was analyzed separately. The results for chromosome arm 3R, which are typical, are shown (Figure 2C-D). From inspection of the fraction of intervals with no insertions (Figure 2C) and from the number of intervals with more insertions than expected by chance alone (Figure 2D), it is clear that the three transposons interact distinctively with the genome.

Minos (Figure 2C-D, red) closely approximates a random distribution. Only 15% more genomic intervals lacked an insert than expected for perfectly random integration, and only a small number of weak candidate hotspots showed up as an excess of intervals with more insertions than expected. An interval could contain more insertions than average due to the presence of a single hotspot, several weaker hotspots, or many dispersed insertions. Candidate *Minos* “hotspots” were usually broader than a single gene. No relationship could be found

between the genes located in different MB hotspots (Table S1). The most striking one was located within a cluster of 25 genes encoding CHK-like kinases (Figure 2E). On either side of this cluster, the density of MB insertions returned to normal.

Both *piggyBac* (Figure 2C, blue) and *P* (Figure 2C, purple) elements showed much greater departures from random integration. Non-random *piggyBac* site-specificity caused the number of unhit intervals to increase by about 30%, whereas *P* insertions left more than twice as many intervals unhit than expected by chance. Both elements also showed a very large excess of hotspots, both in number and in hit frequency (Figure 2D). *P*-element hotspots have been analyzed previously (BELLEN et al. 2004), but it is still unknown why they are targets for preferential insertion. Interestingly, the strongest *piggyBac* hotspot genes (Table S2) significantly differ from those preferentially targeted by *P* elements (Bellen et al. 2004). *piggyBac* target genes frequently encode transcription factors, chromatin factors, and genes involved in growth, nervous system development and behavior.

Differences in transposition relative to genes: Comparing the location of insertions relative to annotated transcripts revealed additional aspects of how these elements target the genome (Figure 3). Intergenic insertions were defined as those lying outside the transcription unit and its promoter, which was assumed to extend 500 bp 5' to the annotated transcription start. *Minos* transposed into such regions 36% of the time, more frequently than either *P* elements (12%) or *piggyBac* elements (24%). The low frequency of *P* element insertion within intergenic regions may result from the strong proclivity of these elements to insert near promoters. About 73% of *P* element insertions (83% of insertions in annotated genes) lie within 500 bp of an annotated 5' transcription start site. In contrast, only 30% of *piggyBac* and 9% of MB insertions were in promoters by this definition. Each time the annotation is revised new promoters are

mapped to more of the orphan *P* insertion sites.

One important potential use of transposon insertions is to generate new gene trap alleles (MORIN *et al.* 2001; BUSZCZAK *et al.* 2007; QUIÑONES-COELLO *et al.* 2007). Gene fusions to GFP (gene traps) can be produced *in vivo* by the transposition of an element bearing splice donor and acceptor sites flanking a GFP-encoding exon. To generate a productive fusion, it is necessary that the transposon integrate into a coding intron of the appropriate splice frame and orientation. Despite the fact that 36% of MB elements insert outside of transcription units, MB elements produced the highest frequency of transposition into coding introns among the three elements tested (Figure 3). MB elements were much better than *P* elements (which were hampered by their promoter bias) but only slightly better than *piggyBac*.

Genome tiling for local mutagenesis: The ability of a transposon to integrate broadly and in effect to tile the genome is critically important for the insertions to be used to manipulate the surrounding region of the genome. To assess breadth of coverage, we plotted the fraction of 40 kb intervals hit at least once within chromosome 3R, as an example, as a function of lambda (λ), the ratio of number of insertions divided by the number of intervals (Figure 4A). At $\lambda = 3$, approximately 95% of intervals will be hit by random insertion (yellow), and our experiments show that *Minos* elements (red) hit about 90%. In contrast, the same number of *P* element insertions (purple) hit only 55% of intervals and *piggyBac* insertions (blue) hit only 77%. How these curves approach saturation will be discussed below, but Figure 4A makes clear that the genome could be quite thoroughly tiled by generating a collection of *Minos* elements equivalent in size or only slightly larger than the current MB collection ($\lambda = 10,458 \times 40 \text{ kb} / 117,000 \text{ kb} = 3.8$).

Polycomb-regulated regions correspond to transposon coldspots: To determine

whether the MB curve in Figure 4A will eventually reach 100% or if there are intervals that cannot be hit by *Minos* insertions, we investigated all 30 examples where two or more adjacent 40 kb zones lacked any insertions at $\lambda = 4$ (excluding 10 basal chromosome regions whose high repetitive DNA content probably impeded mapping). The 30 double negative regions were strikingly non-random and suggested a biological mechanism limiting *Minos* insertion (Table S3). The two largest transposon-free zones occurred on chromosome 3R and corresponded precisely with the BX-C (Figure 4B) and ANT-C homeotic gene complexes. These complexes are known to be regulated at the level of chromatin structure by Polycomb and Trithorax group genes (RINGROSE and PARO 2007). The failure to recover insertions in these domains is not serendipitous, as 17 other MB-free sites also correspond to Polycomb Group (PcG)-regulated gene clusters, including domains that house *ct*, *ems*, *trh*, *nub*, *esg*, *Vsx-1*, *Lim1*, *disco*, and *OdsH*. Direct inspection showed that many other such regions smaller than 80 kb, which would not have been flagged in our analysis, also contained few if any MB insertions. However, not all PcG targets were coldspots; for example, there were many MB insertions in *bru-3* (Figure 4C).

To investigate whether these PcG-regulated domains are coldspots for transposon insertion generally, we also examined whether *P* or *piggyBac* elements integrate normally into these same regions. As can be seen in the case of BX-C (Figure 4B), transposition of *piggyBac* and *P* elements is reduced in PcG-regulated domains as well (Figure 4B). However, some loci appeared to suppress transposon insertion selectively. For example, the region surrounding the PcG-regulated *esg* gene lacked *Minos* inserts, but contained many *piggyBac* and *P* element insertions; indeed, *esg* is a *P*-element hotspot (Figure 4D). Interestingly, *piggyBac* insertions within many such PcG-target regions, including *bru-3* and the *esg* region, were largely “f” class elements (*PBac{WH}*, Table 1, Table S3), suggesting that the engineered structure of the

construct and not just the transposon type affects transposition or marker gene expression within such domains.

Coldspots for *piggyBac* frequently encode membrane proteins: We carried out a similar analysis of *piggyBac* insertions (Figure S1) and identified coldspots that account at least partly for the slower saturation curve of *piggyBac* relative to random integration or *Minos* integration (Table S4). Some were in PcG-target genes that corresponded to sites with reduced MB insertion, although the coldspots were not identical for the two elements. Most interesting, however, was a new class of sites that display normal levels of MB insertion, but exhibit strongly reduced levels of *piggyBac* insertion. These domains are not PcG targets, but are highly enriched in a class of genes with seemingly related function. For example, coldspots include clustered genes encoding acetylcholine receptors (*nAcR α -96* (Figure S1A), *nAcR α -7E*), olfactory or gustatory receptors (*Or69A*, *Or92A*, *Or98A*, *Or22c*, *Gr36a-d*), neuropeptide receptors (*DmsR*, *dpr10*, *CG10418*), GRHRII receptor (*GRHRII*), receptor protein tyrosine phosphatases (*Lar* (Figure S1B), *PTP99A*), dopamine receptors (*D2R*) and ryanodine receptor (*Rya-r44F* (Figure S1C)). Many of the genes encode other putative membrane proteins, often members of the Ig superfamily (*Beat-IIIa*, *Beat-Vc*, *Dpr*, *Dpr2*, *Dpr3*, *Dpr5*, *sns* (Figure S1C)) and channels/transporters (*Glut1*, *Rh50*, *Oatp58Da-c* cluster, *Ir11a*). We conclude that a group of genes with roles in neuronal function, signaling and growth are coldspots for insertion by *piggyBac* elements.

Coldspots for *P* elements include many clustered specialized genes: *P* elements were absent from most of the PcG targets that were also low in MB or *piggyBac* insertion, including ANT-C and BX-C (Figure 3B). Some, but not all, of the domains refractory to *piggyBac* insertion were also low in *P*-element insertions (e.g. Figure S1, Table S4). However, the most

frequent and quantitatively significant classes were intervals containing clusters of genes that were not targeted by *P* elements, but were hit by one or both of the other two transposons. For example, the 20-gene Osiris family (DORER *et al.* 2003) represents one such cluster (Figure S2A). Other clusters, such as the 11-gene esterase complex in region 84D, are mostly refractory to insertion by *P* elements, except for *alpha-Est10*, which was hit 45 times (Figure S2B). In contrast, the 15 MB insertions and 10 *piggyBac* insertions in this region were spread more widely. The MB element in particular was able to insert in many clusters seemingly refractory to *P* element insertion. For example, the MB screen included multiple insertions in eggshell protein genes, genes that have never been hit by *P* elements (Figure S2C).

DISCUSSION

The current GDP collection: The GDP has now generated tools to help functionally analyze at least 9,440 genes, approximately two-thirds of all annotated *Drosophila* proteinencoding genes (Figure 1). Achieving this level of saturation using forward insertional mutagenesis required three different transposons and more than 200,000 independent transpositions. At the time the project began, the genome was not sequenced and relatively little was known about the physical organization of fly genes and regulatory elements. As the project progressed, *Drosophila* researchers and the fly genomics community increasingly documented multiple transcript isoforms, novel RNA genes and key genomic regulatory elements. In response, the GDP project evolved beyond the concept of one disruption per gene, and now comprises more than 14,000 strains. New lines provide additional value by disrupting specific promoters or isoforms, and by providing access to unannotated genes, putative regulatory sequences, and still unknown aspects of genome function.

Recombination-based strategy for completing the GDP: The results reported here

make clear that it would be extremely difficult to achieve 95% genome saturation by random insertional mutagenesis with *P* and *piggyBac* elements alone. Switching to the *Minos* element increased the yield of novel gene hits, but achieving 95% saturation of genes by *Minos* transposition would require an impractically large number of additional insertions to be generated and screened. Including attP sites in a Minos transposon will greatly enhance the general usefulness of insertions for manipulating the genome, since any DNA of interest could be subsequently added at the site of integration. Incorporating DNA that disrupts local chromatin structures might mutate nearby genes, but this approach would be similar to generating deletions from current insertions by imprecise excision. Homologous recombination (RONG et al. 2002) would provide the most attractive finishing strategy for the project, but it has not been technically and economically feasible to carry out on a large enough scale.

Our results suggest a hybrid strategy based on attP-containing Minos insertions for providing access to the remaining genes. An attP site located near a target gene allows efficient homologous recombination by the SIRT method (GAO et al. 2008; GAO et al. 2009). A local duplication containing the mutation of interest is inserted at the attP site and then resolved by generating a local double-strand break (GAO et al. 2008; 2009; reviewed in WESOLOWSKA AND RONG, 2010). Our results shows that *Minos* could be used to tile the entire genome with *attP*-bearing insertions approximately every 40 kb, allowing the efficient application of homologous recombination to disrupt remaining unhit genes. Generating such a collection of elements represents a highly attractive finishing strategy for the GDP and would provide powerful framework for future *Drosophila* genetic manipulation.

A dataset for deducing transposon-genome interactions: A further contribution of the GDP is the detailed knowledge it provides on how transposons interact with their genome. Many

previous studies have demonstrated that specific transposons show a wide variety of non-random integration preferences (reviewed in WU and BURGESS 2004). Some elements are constrained to strongly preferred or invariant target sites by encoded nucleases; for example, *piggyBac* elements only insert at TTAA and *Minos* at TA motifs. In addition, chromatin structure further biases the spectrum of recovered insertion sites (BABENKO et al. 2010; BELLEN *et al.* 2004; GALVAN *et al.* 2007; GANGADHARAN *et al.* 2010; GRABUNDZIJA *et al.* 2010; SIMONS *et al.* 2006; ZHANG and SPRADLING 1994; WALLRATH AND ELGIN, 1995; YAN et al. 2002). However, in metazoans it has usually been difficult to separate site preferences from biases introduced by experimental design, by the loss of marker expression following insertion in suppressive chromatin, and by the failure to accurately map insertions in repetitive DNA.

The GDP datasets of MB and EY transpositions were largely free of bias, as every transposition event from the starting chromosomes that supports marker expression was recovered and analyzed. Quality flanking sequence data were obtained from both the 5' and 3' ends of most insertions, and automated alignments that failed to localize insertions uniquely were usually checked by a human annotator and frequently could be successfully mapped even within repeat-rich genomic regions. The number of insertions with repetitive flanking sequences that could not be mapped uniquely was relatively small (3-5% of total) and consisted of insertions within euchromatic transposons or within repetitive segments of centric heterochromatin and the Y chromosome. Thus, with respect to potential bias from both chromatin and repetitive genomic sequences, GDP data provide an accurate picture of transposon site-selectivity within euchromatin, but an incomplete picture of transposition within centric heterochromatin.

Transposons avoid insertion within many Polycomb-regulated regions: Our data show how chromatin structure influences transposon insertion. In particular, many regions in the

genome enriched in the repressive histone modification H3K27me3 were targeted much less frequently by all three transposons. Repressive domains frequently arise from the activity of Polycomb group genes (SCHWARTZ *et al.* 2006; 2010). Many such regions contain clustered genes encoding key transcription factors such as Hox genes that regulate tissue differentiation and development. Each such cluster is repressed in some cells during development but active in others. Consequently, the roster of PcG-repressed domains depends on the cell type in question. In yeast, plants and *Drosophila*, H3K27me-rich centric heterochromatin is likely to be generated using other pathways, including the piRNA pathway (reviewed in RIDDLE AND ELGIN, 2008). Our studies focused on germline transposition, which cluster size analysis places during premeiotic and meiotic adult germ cell development.

The observation that transposon insertions are recovered less frequently in PcG-regulated, “closed” domains has been reported previously (BELLEN *et al.* 2004; GRABUNDZIJA *et al.* 2010; SIMONS *et al.* 2006). For example, *Tol2* integrations are underrepresented in regions rich in H3K27me in human cells (GRABUNDZIJA *et al.* 2010), however, *piggyBac* insertions are not. In most cases it was difficult to determine if the dearth of insertions was due to blocked transposition or reduced marker expression, however.

Our data suggest that PcG-regulated regions directly suppress transposition, but they likely also reduce marker gene expression. The *yellow* gene and the *Pax6*-GFP construct used to detect EY and MB transpositions are sufficiently robust to detect at least some insertions in centric heterochromatin. Hence these elements must actually transpose with reduced frequency into PcG-domains because most such insertions would be detected. Consistent with this view, when suppressors of variegation were used to reveal the location of “suppressed” insertions they were only found in centric heterochromatin (ZHANG and SPRADLING 1994; YAN *et al.* 2002). A

direct effect on transposition is less certain in the case of *piggyBac*, because the elements studied carry the position-effect sensitive *mini-white* gene, and “f” insertions, which carry a chromatin insulator, were preferentially recovered in some PcG-domains (Table S3). Our results suggest that as in pluripotent mammalian cells (BOYER ET AL. 2007), *Drosophila* PcG-domains are already established in pre-meiotic germ cells, where they can affect adult germline transposition. Functions of Polycomb genes in the early germline have been described in the male (Chen et al. 2005).

Transposon-free genomic regions: The genomes of many organisms, including humans, contain rare “transposon-free regions (TFRs).” Some of these encode clustered HOX genes such as the human *HOXA4-11*, *HOXB4-6* and *HOXD8-13* loci (SIMONS et al. 2006) that resemble the *Drosophila* BX-C and ANT-C complexes, which we observed to resist transposon insertion. We looked to see if other human TFRs have *Drosophila* homologs that are also refractory to integration. For example, human *DLX5* lies within a TFR, and its *Drosophila* homolog *Distalless* is a PcG-regulated gene that was not hit by any of the three transposons. Many TFRs did not show such a correlation, however. *PAX6* lies within a TFR and the closely related *Drosophila* genes *eyeless (ey)* and *sine oculis (so)* both lie within PcG-regulated domains (SCHWARTZ et al. 2010). However, *ey* received one *piggyBac* and two MB insertions in our experiments, while the *so* region was hit by two MB insertions. Finally, the region surrounding the *NR2F1/COUP-TF1* gene is a TFR in at least six vertebrate genomes (SIMONS et al. 2006). *Drosophila seven-up* is a *COUP-TF1* homolog, however, it does not lie within a PcG-regulated domain and was the target of ten MB and three *piggyBac* insertions in our experiments. The domains that were refractory to insertion in our experiments frequently contain natural integrated

transposons in some strains. Thus, in *Drosophila* suppression of transposon activity by PcG-regulated domains appears insufficient to sustain transposon-free regions. If PcG-mediated repression of transposon insertion is important in the genesis of mammalian TFRs, it may exert stronger effects, synergize with other regulatory mechanisms not present in *Drosophila*, and act on rates of germline transposition that are much lower than those in *Drosophila*.

Some transposons may evolve to benefit their host: Transposon insertions frequently disrupt vital genes, hence the introduction and spread of transposable elements within a genome has the potential to be highly deleterious. Consequently, like viruses, transposons should evolve to minimize costs to host fitness. In addition, increasing evidence documents a major creative role for transposable elements in the evolution of new genes, regulatory elements and on genome size itself (SINZELLE et al. 2009). A transposon that could generate useful variation within the genome of its host under conditions of stress, might contribute to the survival of both its host and itself (MCCLINTOCK, 1984). An element might minimize damage and maximize the chance of adaptive variation by avoiding insertion in evolutionarily stable genes, and selectively targeting genes whose structure and/or regulation evolves rapidly.

We observed several examples of site-specificity in our experiments that suggested such an adaptation. The gene cluster encoding proteins with CHK-like kinase domains (Figure 2E) was one of few hotspots for *Minos* insertion. One of these genes, *CHKov1*, has been shown to harbor a *Doc* insertion in many wild *Drosophila* populations that confers enhanced insecticide resistance (AMINETZACH et al. 2005). *piggyBac* elements rarely inserted in genes that encode a variety of membrane receptors for neurotransmitters and other ligands (Table S4). Conversely, many *piggyBac* hotspot loci contain genes affecting neural development and behavior (Table S2). Thus, of the three elements studied, *piggyBac* was the only one whose site preferences were

suggestive of having evolved to minimize damage and to maximize changes in the regulation of potentially adaptive genes following insertion. There may be common transcription factors or chromatin configurations at these sites that allow such targeting.

Implications for other organisms utilizing transposon mutagenesis: Recently there has been growing interest in the application of insertional mutagenesis in a wide variety of experimental organisms both in the germ line (DING et al. 2005, BAZOPOULOU AND TAVERNARAKIS, 2009; O'MALLEY AND ECKER, 2009; GALVÁN et al. 2007; SIVASUBBU et al. 2007; SASAKURA et al. 2007; DE WIT et al. 2010) and in somatic cells (review: COPELAND AND JENKINS, 2010). Indeed, the potential application of transposons as human gene therapy vectors is currently undergoing clinical trials (IZSVÁK et al. 2010). The lessons learned in the GDP project regarding both the common and unique ways that transposons interact and evolve with the genome are certain to help these projects maximize the value of these exceptional tools for natural and human-guided genetic manipulation.

ACKNOWLEDGMENTS

We thank Danqing Bei, Ying Fang, Adeel Jawaaid, Jianping Li, Zhihua Wang, and Jin Yue for generating and maintaining fly stocks. Vanessa Damm, Shelly Paterno and Eric Chen assisted in the line maintenance and balancing at Carnegie. We thank Soo Park and Kenneth Wan at LBNL for assistance with iPCR and sequencing of insertions. We are grateful to Exelixis, Inc. and Aprogen, Inc. (formerly GenExel) for providing lines and sequence data. We are grateful to researchers at Max Planck Göttingen, EMBL Heidelberg and at DeveloGen AG for donating *P* insertion lines to the public. We are grateful to Ulrich Schäfer and Herbert Jäckle for providing information and lines from the Göttingen X-linked insertion collection. We thank Peter Maroy for shipping lines to the project from the Szeged stock center. We thank Kathy Matthews, Kevin Cook and Annette Parks for coordinating the transition of the lines to the BDSC. We thank Koen Venken for useful suggestions. This work was supported by NIGMS (GM067858). Additional funds were provided through the support of the Spradling and Bellen labs from the Howard Hughes Medical Institute.

REFERENCES

- ADAMS, M. D., S. E. CELNIKER, R. A. HOLT, C. A. EVANS, J. D. GOCAYNE *et al.*, 2000 The genome sequence of *Drosophila melanogaster*. *Science* **287**: 2185-2195.
- AMINETZACH, Y. T., J. M. MACPHERSON and D. A. PETROV, 2005 Pesticide resistance via transposition-mediated adaptive gene truncation in *Drosophila*. *Science* **309**: 764-767.
- BABENKO, V. N., I. V. MAKUNIN, I. V. BRUSENTOVA, E. S. BELYAEVA, D. A. MAKSIMOV *et al.*, 2010 Paucity and preferential suppression of transgenes in late replication domains of the *D. melanogaster* genome. *BMC Genomics* **11**: 318.
- BATEMAN J.R., A.M. LEE, AND C.T.WU, 2006 Site-specific transformation of *Drosophila* via phiC31 integrase-mediated cassette exchange. *Genetics* **173**: 769-77.
- BAUDRY, C., S. MALINSKY, M. RESTITUITO, A. KAPUSTA, S. ROSA *et al.*, 2009 PiggyMac, a domesticated piggyBac transposase involved in programmed genome rearrangements in the ciliate *Paramecium tetraurelia*. *Genes Dev* **23**: 2478-2483.
- BAZOPOULOU, D., and N. TAVERNARAKIS, 2009 The NemaGENETAG initiative: large scale transposon insertion gene-tagging in *Caenorhabditis elegans*. *Genetica* **137**: 39-46.
- BEINERT, N., M. WERNER, G. DOWE, H.-R. CHUNG, H. JÄCKLE *et al.*, 2004 Systematic gene targeting on the X chromosome of *Drosophila melanogaster*. *Chromosoma* **113**: 271-275.
- BELLEN, H. J., R. W. LEVIS, G. LIAO, Y. HE, J. W. CARLSON *et al.*, 2004 The BDGP gene disruption project: single transposon insertions associated with 40% of *Drosophila* genes. *Genetics* **167**: 761-781.
- BELLEN, H. J., C. J. O'KANE, C. WILSON, U. GROSSNIKLAS, R. K. PEARSON *et al.*, 1989 *P*-element-mediated enhancer detection: a versatile method to study development in *Drosophila*. *Genes Dev* **3**: 1288-1300.

- BENDER, W., and A. HUDSON, 2000 P element homing to the *Drosophila* bithorax complex. *Development* **127**: 3981-3992.
- BIER, E., H. VAESSIN, S. SHEPHERD, K. LEE, K. MCCALL *et al.*, 1989 Searching for pattern and mutation in the *Drosophila* genome with a *P-lacZ* vector. *Genes Dev* **3**: 1273-1287.
- BOYER, L. A., K. PLATH, J. ZEITLINGER, T. BRAMBRINK, L. A. MEDEIROS *et al.*, 2006 Polycomb complexes repress developmental regulators in murine embryonic stem cells. *Nature* **441**: 349-353.
- BRENNECKE, J., A. A. ARAVIN, A. STARK, M. DUS, M. KELLIS *et al.*, 2007 Discrete small RNA-generating loci as master regulators of transposon activity in *Drosophila*. *Cell* **128**: 1089-1103.
- BRENNECKE, J., D. R. HIPFNER, A. STARK, R. B. RUSSELL and S. M. COHEN, 2003 bantam encodes a developmentally regulated microRNA that controls cell proliferation and regulates the proapoptotic gene *hid* in *Drosophila*. *Cell* **113**: 25-36.
- BUSZCZAK, M., S. PATERNO, D. LIGHTHOUSE, J. BACHMAN, J. PLANCK *et al.*, 2007 The Carnegie protein trap library: a versatile tool for *Drosophila* developmental studies. *Genetics* **175**: 1505-1531.
- CHEN X, M. HILLER, Y. SANCAK AND M.T. FULLER, 2005 Tissue-specific TAFs counteract Polycomb to turn on terminal differentiation. *Science* **310**: 869-72.
- COOLEY, L., R. KELLEY and A. SPRADLING, 1988 Insertional mutagenesis of the *Drosophila* genome with single P elements. *Science* **239**: 1121-1128.
- COPELAND, N. G., and N. A. JENKINS, 2010 Harnessing transposons for cancer gene discovery. *Nat Rev Cancer* **10**: 696-706.

- COUFAL, N. G., J. L. GARCIA-PEREZ, G. E. PENG, G. W. YEO, Y. MU *et al.*, 2009 L1 retrotransposition in human neural progenitor cells. *Nature* **460**: 1127-1131.
- DE WIT, T., S. DEKKER, A. MAAS, G. BREEDVELD, T. A. KNOCH *et al.*, 2010 Tagged mutagenesis by efficient Minos-based germ line transposition. *Mol. Cell Biol.*: **30**: 68-77.
- DING, S., X. WU, G. LI, M. HAN, Y. ZHUANG *et al.*, 2005 Efficient transposition of the piggyback (PB) transposon in mammalian cells and mice. *Cell* **122**: 473-483.
- DORER, D. R., J. A. RUDNICK, E. N. MORIYAMA and A. C. CHRISTENSEN, 2003 A family of genes clustered at the Triplo-lethal locus of *Drosophila melanogaster* has an unusual evolutionary history and significant synteny with *Anopheles gambiae*. *Genetics* **165**: 613-621.
- EWING B, P. GREEN, 1998 Basecalling of automated sequencer traces using phred. II. Error probabilities. *Genome Research* **8**: 186-194 (1998)
- FUJITA, P. A., B. RHEAD, A. S. ZWEIG, A. S. HINRICHS, D. KAROLCHIK *et al.*, 2010 The UCSC Genome Browser database: update 2011. *Nucleic Acids Research* **39**: D876-D882.
- GALVÁN, A., D. GONZÁLEZ-BALLESTER and E. FERNÁNDEZ, 2007 Insertional mutagenesis as a tool to study genes/functions in *Chlamydomonas*. *Adv Exp Med Biol* **616**: 77-89.
- GANGADHARAN, S., L. MULARONI, J. FAIN-THORNTON, S. J. WHEELAN and N. L. CRAIG, 2010 DNA transposon Hermes inserts into DNA in nucleosome-free regions in vivo. *Proc Natl Acad Sci U S A* **107**: 21966-21972.
- GAO, G., C. MCMAHON, J. CHEN and Y. S. RONG, 2008 A powerful method combining homologous recombination and site-specific recombination for targeted mutagenesis in *Drosophila*. *Proc Natl Acad Sci U S A* **105**: 13999-14004.

- GAO, G., N.WESOŁOWSKA and Y. S. RONG, 2009 SIRT combines homologous recombination, site-specific integration, and bacterial recombineering for targeted mutagenesis in *Drosophila*. Cold Spring Harb Protoc **2009**: pdb prot5236.
- GODFREY, A. C., J. M. KUPSCO, B. D. BURCH, R. M. ZIMMERMAN, Z. DOMINSKI *et al.*, 2006 U7 snRNA mutations in *Drosophila* block histone pre-mRNA processing and disrupt oogenesis. Rna **12**: 396-409.
- GRABUNDZIJA, I., M. IRGANG, L.MATES, E. BELAY, J.MATRAI *et al.*, 2010 Comparative analysis of transposable element vector systems in human cells. Mol Ther **18**: 1200-1209.
- GROTH AC, M. FISH, R. NUSSE, M.P. CALOS, 2004 Construction of transgenic *Drosophila* by using the site-specific integrase from phage phiC31. Genetics **166**: 1775-1782.
- HORN, C., B. JAUNICH and E. A.WIMMER, 2000 Highly sensitive, fluorescent transformation marker for *Drosophila* transgenesis. Dev Genes Evol **210**: 623-629.
- IZSVÁK, Z., P. B. HACKETT, L. J. COOPER and Z. IVICS, 2010 Translating Sleeping Beauty transposition into cellular therapies: victories and challenges. Bioessays **32**: 756-767.
- KIM, Y. I., T. RYU, J. LEE, Y. S. HEO, J. AHNN *et al.*, 2010 A genetic screen for modifiers of *Drosophila* caspase Dcp-1 reveals caspase involvement in autophagy and novel caspase-related genes. BMC Cell Biol **11**: 9.
- LIAO, G. C., E. J. REHM and G. M. RUBIN, 2000 Insertion site preferences of the P transposable element in *Drosophila melanogaster*. Proc Natl Acad Sci U S A **97**: 3347-3351.
- METAXAKIS, A., S. OEHLER, A. KLINAKIS and C. SAVAKIS, 2005 *Minos* as a genetic and genomic tool in *Drosophila melanogaster*. Genetics **171**: 571-581.
- MCCLINTOCK B. (1984). The significance of responses of the genome to challenge. Science 226: 792-801.

- MISRA S, M.A. CROSBY, C.J.MUNGALL, B.B.MATTHEWS, K.S. CAMPBELL *et al.*, 2002
Annotation of the *Drosophila melanogaster* euchromatic genome: a systematic review.
Genome Biol. **3**: RESEARCH0083.
- MORIN, X., R. DANEMAN, M. ZAVORTINK and W. CHIA, 2001 A protein trap strategy to detect
GFP-tagged proteins expressed from their endogenous loci in *Drosophila*. Proc Natl Acad
Sci U S A **98**: 15050-15055.
- O'MALLEY, R. C., and J. R. ECKER, 2010 Linking genotype to phenotype using the Arabidopsis
unimutant collection. Plant J **61**: 928-940.
- OH, S.W., T. KINGSLEY, H. H. SHIN, Z. ZHENG, H. W. CHEN *et al.*, 2003 A P-element insertion
screen identified mutations in 455 novel essential genes in *Drosophila*. Genetics **163**:
195-201.
- PETER, A., P. SCHÖTTLER, M.WERNER, N. BEINERT, G. DOWE *et al.*, 2002 Mapping and
identification of essential gene functions on the X chromosome of *Drosophila*. EMBO
Rep **3**: 34-38.
- QUÍÑONES-COELLO, A.T., L.N. PETRELLA, L. AYERS, A.MELILLO, S.MAZZALUPO *et al.*, 2007
Exploring strategies for protein trapping in *Drosophila*. Genetics **175**: 1089-1104.
- RIDDLE, N. C., and S. C. ELGIN, 2008 A role for RNAi in heterochromatin formation in
Drosophila. Curr Top Microbiol Immunol **320**: 185-209.
- RINGROSE, L., and R. PARO, 2007 Polycomb/Trithorax response elements and epigenetic
memory of cell identity. Development **134**: 223-232.
- RONG, Y. S., 2002 Gene targeting by homologous recombination: a powerful addition to the
genetic arsenal for *Drosophila* geneticists. Biochem Biophys Res Commun **297**: 1-5.

- RØRTH, P., 1996 A modular misexpression screen in *Drosophila* detecting tissue-specific phenotypes. *Proc Natl Acad Sci U S A* **93**: 12418-12422.
- ROY, S., J. ERNST, P. V. KHARCHENKO, P. KHERADPOUR, N. NEGRE *et al.*, 2010 Identification of Functional Elements and Regulatory Circuits by *Drosophila* modENCODE. *Science*.
- RYDER, E., F. BLOWS, M. ASHBURNER, R. BAUTISTA-LLACER, D. COULSON *et al.*, 2004 The DrosDel collection: a set of P-element insertions for generating custom chromosomal aberrations in *Drosophila melanogaster*. *Genetics* **167**: 797-813.
- SASAKURA, Y., Y. OOGAI, T. MATSUOKA, N. SATOH and S. AWAZU, 2007 Transposon mediated transgenesis in a marine invertebrate chordate: *Ciona intestinalis*. *Genome Biol* **8 Suppl 1**: S3.
- SCHWARTZ, Y. B., T. G. KAHN, D. A. NIX, X. Y. LI, R. BOURGON *et al.*, 2006 Genome-wide analysis of Polycomb targets in *Drosophila melanogaster*. *Nat Genet* **38**: 700-705.
- SCHWARTZ, Y. B., T. G. KAHN, P. STENBERG, K. OHNO, R. BOURGON *et al.*, 2010 Alternative epigenetic chromatin states of polycomb target genes. *PLoS Genet* **6**: e1000805.
- SIMONS, C., M. PHEASANT, I. V. MAKUNIN and J. S. MATTICK, 2006 Transposon-free regions in mammalian genomes. *Genome Res* **16**: 164-172.
- SINZELLE L, Z. IZSVÁK Z AND Z. IVICS, 2009 Molecular domestication of transposable elements: from detrimental parasites to useful host genes. *Cell Mol Life Sci*. **66**: 1073-93.
- SIVASUBBU, S., D. BALCIUNAS, A. AMSTERDAM and S. C. EKKER, 2007 Insertional mutagenesis strategies in zebrafish. *Genome Biol* **8 Suppl 1**: S9.
- SPRADLING, A. C., D. M. STERN, I. KISS, J. ROOTE, T. LAVERTY *et al.*, 1995 Gene disruptions using P transposable elements: an integral component of the *Drosophila* genome project. *Proc Natl Acad Sci U S A* **92**: 10824-10830.

- SPRADLING, A. C., D. STERN, A. BEATON, E. J. RHEM, T. LAVERTY *et al.*, 1999 The Berkeley Drosophila Genome Project gene disruption project: Single P-element insertions mutating 25% of vital Drosophila genes. *Genetics* **153**: 135-177.
- STAUDT, N., A. MOLITOR, K. SOMOGYI, J. MATA, S. CURADO *et al.*, 2005 Gain-of-function screen for genes that affect Drosophila muscle pattern formation. *PLoS Genet.* 1, e55; 499-506.
- THIBAUT, S. T., M. A. SINGER, W. Y. MIYAZAKI, B. MILASH, N. A. DOMPE *et al.*, 2004 A complementary transposon tool kit for Drosophila melanogaster using P and piggyBac. *Nat Genet* **36**: 283-287.
- TOLHUIS, B., E. DE WIT, I. MUIJRS, H. TEUNISSEN, W. TALHOUT *et al.*, 2006 Genome-wide profiling of PRC1 and PRC2 Polycomb chromatin binding in Drosophila melanogaster. *Nat Genet* **38**: 694-699.
- TOWER, J., and R. KURAPATI, 1994 Preferential transposition of a Drosophila P element to the corresponding region of the homologous chromosome. *Mol Gen Genet* **244**: 484-490.
- TWEEDIE, S., M. ASHBURNER, K. FALLS, P. LEYLAND, P. MCQUILTON *et al.*, 2009 FlyBase: enhancing Drosophila Gene Ontology annotations. *Nucleic Acids Res* **37**: D555-559.
- VENKEN, K. J., Y. HE, R. A. HOSKINS and H. J. BELLEN, 2006 P[acman]: a BAC transgenic platform for targeted insertion of large DNA fragments in D. melanogaster. *Science* **314**: 1747-1751.
- VENKEN, K. J., and H. J. BELLEN, 2007 Transgenesis upgrades for Drosophila melanogaster. *Development* **134**: 3571-3584.
- VENKEN, K. J., J. W. CARLSON, K. L. SCHULZE, H. PAN, Y. HE *et al.*, 2009 Versatile P[acman] BAC libraries for transgenesis studies in Drosophila melanogaster. *Nat Methods* **6**: 431-434.

WALLRATH, L. L., and S. C. ELGIN, 1995 Position effect variegation in *Drosophila* is associated with an altered chromatin structure. *Genes Dev* **9**: 1263-1277.

WESOŁOWSKA N, RONG YS. (2010). The past, present and future of gene targeting in *Drosophila*. *Fly* 4: 53-9.

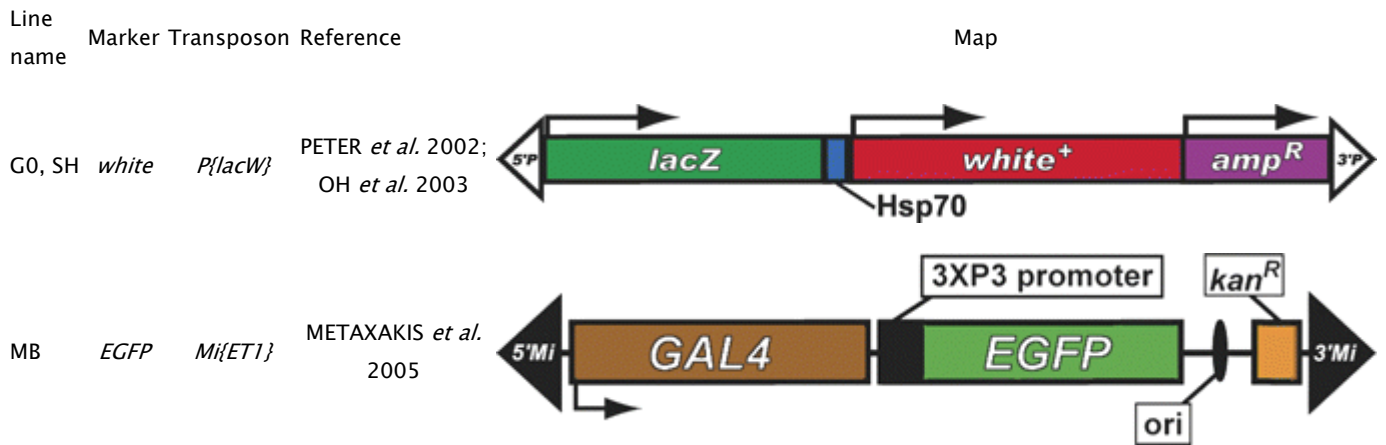
WU, X., and S. M. BURGESS, 2004 Integration target site selection for retroviruses and transposable elements. *Cellular and Molecular Life Sciences* **61**: 2588-2596.

YAN, C. M., K. W. DOBIE, H. D. LE, A. Y. KONEV and G. H. KARPEN, 2002 Efficient recovery of centric heterochromatin P-element insertions in *Drosophila melanogaster*. *Genetics* **161**: 217-229.

ZHANG, P., and A. C. SPRADLING, 1994 Insertional mutagenesis of *Drosophila* heterochromatin with single P elements. *Proc Natl Acad Sci U S A* **91**: 3539-3543.

Table 1
Mutator transposons

Line name	Marker	Transposon	Reference	Map
EY	<i>white</i> , <i>yellow</i>	<i>P{EPgy2}</i>	BELLEN <i>et al.</i> 2004	
HP	<i>white</i>	<i>P{EPg}</i>	STAUDT <i>et al.</i> 2005	
DP, GG	<i>yellow</i>	<i>P{Mae-UAS.6.11}</i>	Beinert <i>et al.</i> 2004; STAUDT <i>et al.</i> 2005	
d	<i>white</i>	<i>P{XP}</i>	THIBAUT <i>et al.</i> 2004	
c	<i>white</i>	<i>PBac{PB}</i>	THIBAUT <i>et al.</i> 2004	
e	<i>white</i>	<i>PBac{RB}</i>	THIBAUT <i>et al.</i> 2004	
f	<i>white</i>	<i>PBac{WH}</i>	THIBAUT <i>et al.</i> 2004	
G	<i>white</i>	<i>P{EP}</i>	RØRTH 1996; KIM <i>et al.</i> 2010; GenExel Library at KAIST (http://genexel.kaist.ac.kr/mapvie_w3/index.html)	



*The schematic diagrams are not drawn to scale and are meant only to indicate the components present in each transposon. Thin lines separating some components have been added to prevent labels from overlapping and are not intended to indicate spacers between components. Please refer to the original publications and curated FlyBase reports for details.

Table 2
Summary of GDP lines

Collection	BDSC Lines	In genes	Intergenic	New genes
SPRADLING <i>et al.</i> 1999	934	898	36	936
BELLEN <i>et al.</i> 2004	6062	5118	944	3910
New EYs	1193	1059	134	641
Exelixis	1859	1800	59	1983 ^a
STAUDT <i>et al.</i> 2005 ^b	284	276	8	109
MB	2658	2147	511	1155
GenExel	1136	1120	16	616
Other	530	514	16	90
Total	14,656	12,932	1724	9440

The numbers of strains from the indicated sources selected for the GDP collection and currently available at the BDSC are shown. The numbers of strains containing insertions in genes (see *Methods*) or within intergenic regions are also given. The *New Genes* column gives the number of genes hit by insertions in that collection that are not hit by insertions from the collections above it in the table. The values reflect the current status of the GDP collection; the values for the SPRADLING *et al.* 1999 and BELLEN *et al.* 2004 collections are lower than those originally reported, due to loss or replacement with strains hitting the same gene from later collections (see *Methods*).

^a Includes 357 genes hit by lines that were sent to the Harvard Stock Center, rather than BDSC.

^b The STAUDT *et al.* 2005 collection is also referred to as the Max Plank/EMBL/DeveloGen collection in *Methods*.

FIGURE LEGENDS

FIGURE 1. Growth of the GDP strain collection

The total number of GDP strains (green triangles) and the number of genes with one or more associated GDP lines (filled circles) are shown as a function of time beginning with the completion of the *Drosophila* genome sequence in 2000, which signaled the end of project phase 1. In 2010, the project completed phase 2 in which genes were targeted based on the location of insertions from undirected forward screens.

FIGURE 2. Saturation behavior of *P*, *piggyBac* and *Minos* insertions

(A) Plot of MB insertions per 250 kb vs interval number along chromosome 3L reveals a large hotspot. (B) MB insertions within 10 kb intervals around the hotspot in A. The number per interval expected by chance is shown in pink. 0 corresponds to 3L:12580233, the site on the homolog of the mobilized element in the MB screen. (C-D) Distribution of MB (red), *piggyBac* (blue) or EY (purple) insertions within 10 kb genomic intervals on chromosome 3R, compared with random transposition (Poisson distribution, yellow). To facilitate comparison, the same numbers of insertions were analyzed in each case (2790; corresponding to 1 insertion per interval). The number of intervals with 0 insertions (C, "0") is relevant to coldspot behavior; intervals hit more frequently than by random expectation (D) are indicative of *piggyBac* and *P* element hotspots. (E) The *Minos* hotspot located within a cluster of genes encoding CHK-kinases on chromosome 3R. The locations of MB (*Minos*), Pig (*piggyBac*), and EY (*P*) element insertions are shown by vertical bars above the gene map of the region.

FIGURE 3. Transposon insertion with respect to transcript structure.

The percentage of MB, *piggyBac* (Pig) and EY insertions located in the indicated regions of annotated transcripts are shown. Numbers may not sum to 100% because an insertion may disrupt multiple transcripts in different positions. A region was scored positive if one or more annotated transcripts with the indicated character were hit by an insertion. To simplify calculation, only the first four annotated transcripts hit by the insertion were considered in determining these values. Because of the large N values, the 95% confidence intervals of these proportions were always less than +/-1%. Consequently,

the differences were significant except in the case of MB compared to EY insertion in non-coding introns.

FIGURE 4. Transposons non-randomly avoid some genomic intervals, including regions with PcG-dependent repressive marks.

(A) The saturation behavior of 40 kb genomic intervals for transposon insertion on chromosome 3R is plotted as λ (the ratio of the number of insertions/ the number of intervals) increases. Poisson (random) expectation (yellow), MB (*Minos*) elements (red), *piggyBac* elements (blue), EY (*P*) elements (purple). EY elements saturate well below 100%. In contrast, MB elements approach saturation only slightly more slowly than random, whereas *piggyBacs* appear intermediate. (B) MB, *piggyBac* and EY elements insert with greatly reduced frequency in the *Ultrabithorax* gene cluster. Regions of the *Drosophila* genome as displayed on the UCSC browser are shown. Insertion sites for these elements are shown in labeled tracks above the map as vertical lines of unit thickness (MB in red; *piggyBac* in blue; EY in purple; thicker lines denote multiple insertions). The orange boxes display the approximate position of PcG target regions as mapped by Schwartz et al. (2010). (C) Similar display of the *bru-3* gene region shows that not all Polycomb-regulated chromatin domains are transposon-poor. (D) The *esg* gene cluster and its surrounding region illustrates that some PcG targets are largely refractory to MB insertion, but not to the other two elements.

Figure 1

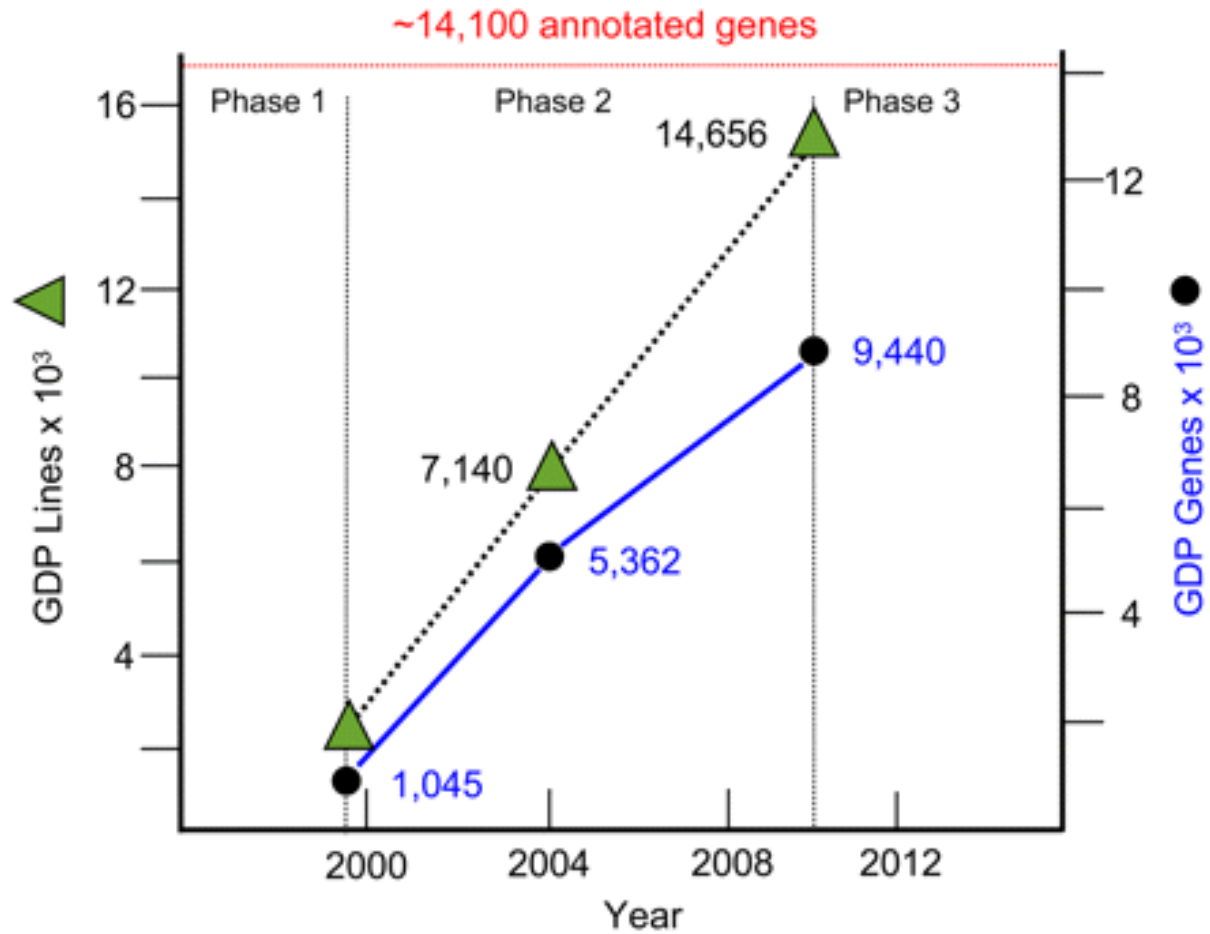


FIGURE 1. Growth of the GDP strain collection

The total number of GDP strains (green triangles) and the number of genes with one or more associated GDP lines (filled circles) are shown as a function of time beginning with the completion of the *Drosophila* genome sequence in 2000, which signaled the end of project phase 1. In 2010, the project completed phase 2 in which genes were targeted based on the location of insertions from undirected forward screens.

Figure 2

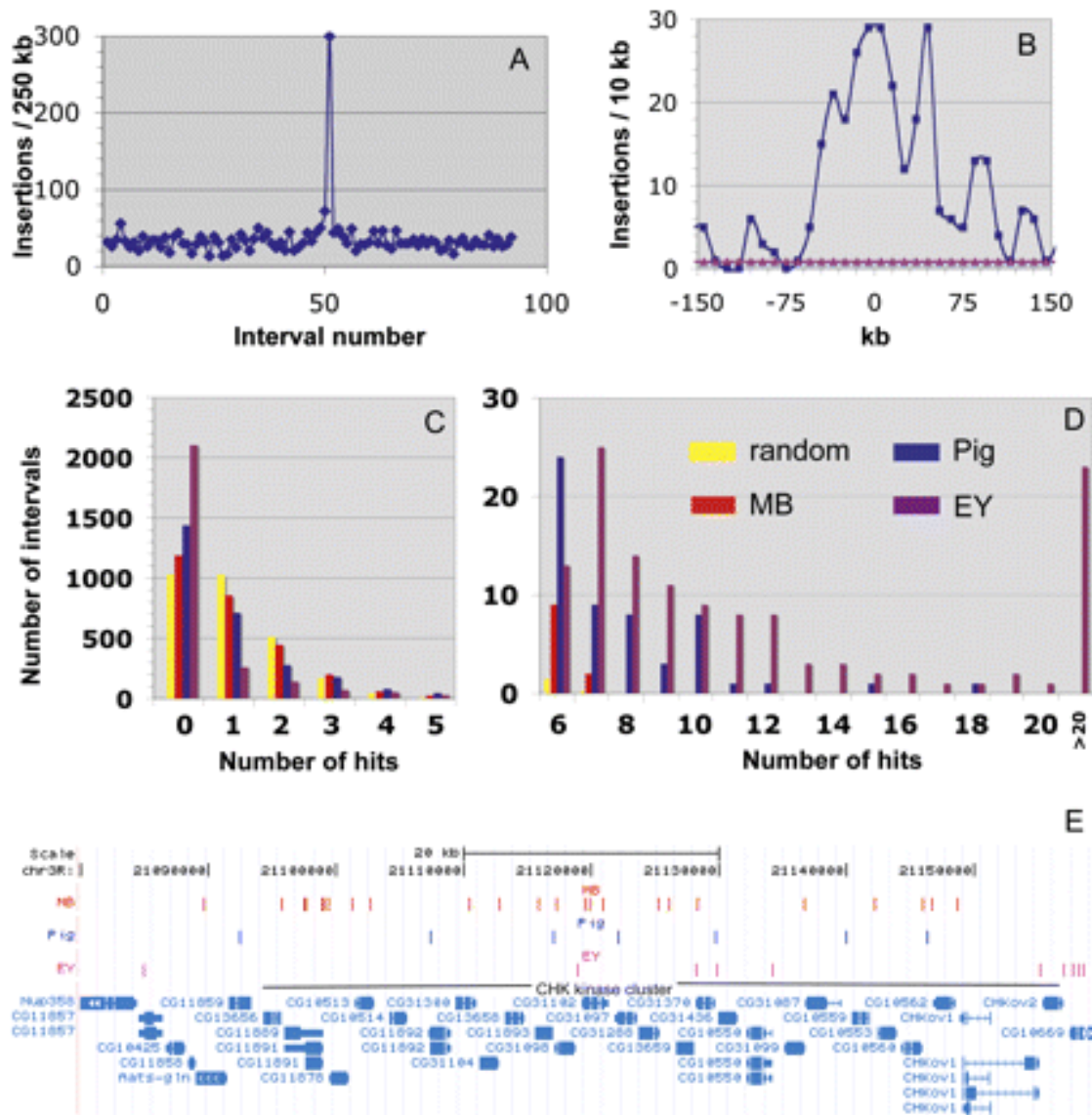


FIGURE 2. Saturation behavior of *P*, *piggyBac* and *Minos* insertions

(A) Plot of MB insertions per 250 kb vs interval number along chromosome 3L reveals a large hotspot. (B) MB insertions within 10 kb intervals around the hotspot in A. The number per interval expected by chance is shown in pink. 0 corresponds to 3L:12580233, the site on the homolog of the mobilized element in the MB screen. (C-D) Distribution of MB (red), *piggyBac* (blue) or EY (purple) insertions within 10 kb genomic intervals on chromosome 3R, compared with random transposition (Poisson distribution, yellow). To facilitate comparison, the same numbers of insertions were analyzed in each case (2790; corresponding to 1 insertion per interval). The number of intervals with 0 insertions (C, "0") is relevant to coldspot behavior; intervals hit more frequently than by random expectation (D) are indicative of *piggyBac* and *P* element hotspots. (E) The *Minos* hotspot located within a cluster of genes encoding CHK-kinases on chromosome 3R. The locations of MB (*Minos*), Pig (*piggyBac*), and EY (*P*) element insertions are shown by vertical bars above the gene map of the region.

Figure 3

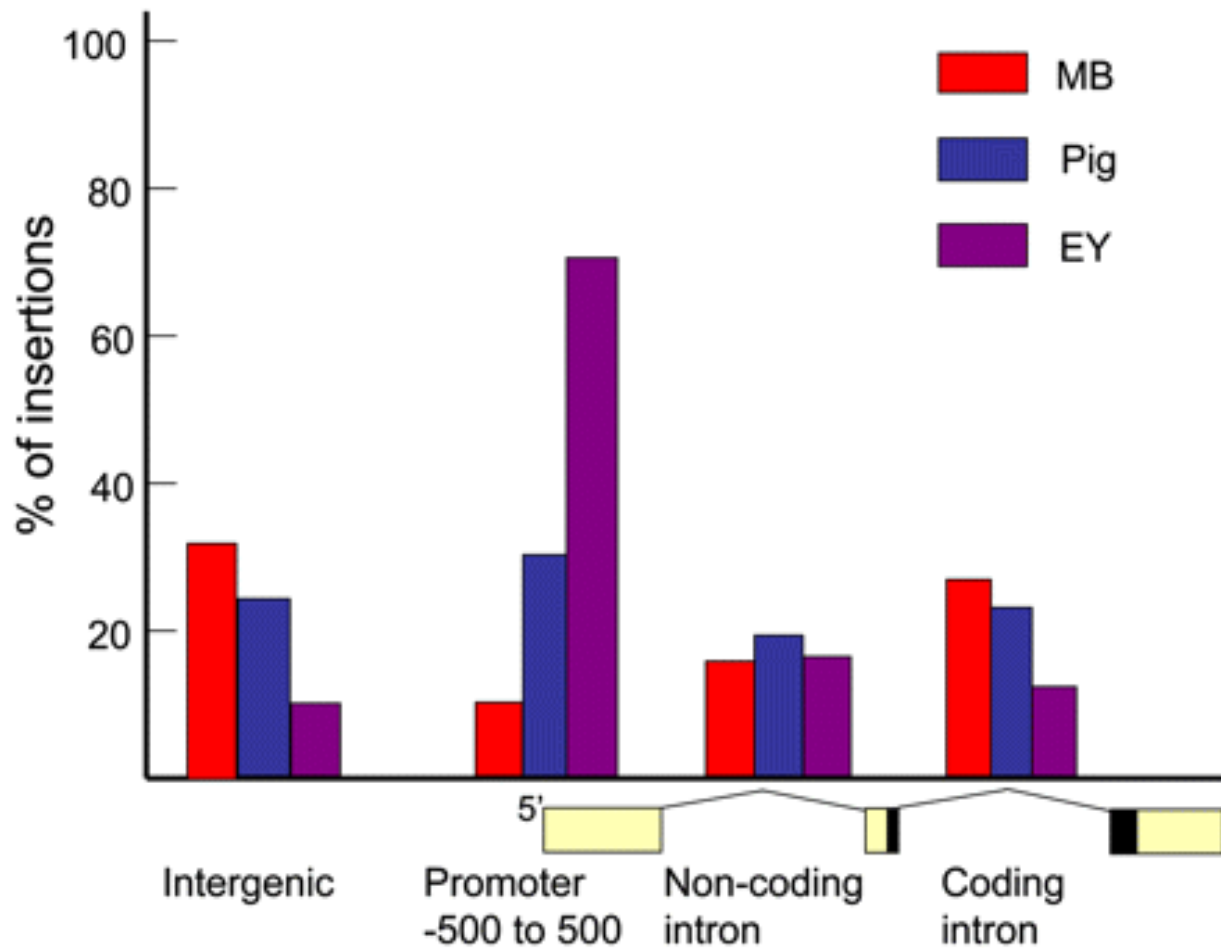


FIGURE 3. Transposon insertion with respect to transcript structure.

The percentage of MB, *piggyBac* (Pig) and EY insertions located in the indicated regions of annotated transcripts are shown. Numbers may not sum to 100% because an insertion may disrupt multiple transcripts in different positions. A region was scored positive if one or more annotated transcripts with the indicated character were hit by an insertion. To simplify calculation, only the first four annotated transcripts hit by the insertion were considered in determining these values. Because of the large N values, the 95% confidence intervals of these proportions were always less than $\pm 1\%$. Consequently, the differences were significant except in the case of MB compared to EY insertion in non-coding introns.

(A) The saturation behavior of 40 kb genomic intervals for transposon insertion on chromosome 3R is plotted as λ (the ratio of the number of insertions/ the number of intervals) increases. Poisson (random) expectation (yellow), MB (*Minos*) elements (red), *piggyBac* elements (blue), EY (*P*) elements (purple). EY elements saturate well below 100%. In contrast, MB elements approach saturation only slightly more slowly than random, whereas *piggyBacs* appear intermediate. (B) MB, *piggyBac* and EY elements insert with greatly reduced frequency in the *Ultrabithorax* gene cluster. Regions of the *Drosophila* genome as displayed on the UCSC browser are shown. Insertion sites for these elements are shown in labeled tracks above the map as vertical lines of unit thickness (MB in red; *piggyBac* in blue; EY in purple; thicker lines denote multiple insertions). The orange boxes display the approximate position of PcG target regions as mapped by Schwartz et al. (2010). (C) Similar display of the *bru-3* gene region shows that not all Polycomb-regulated chromatin domains are transposon-poor. (D) The *esg* gene cluster and its surrounding region illustrates that some PcG targets are largely refractory to MB insertion, but not to the other two elements.